

Rec'd PCT/PTO 04 OCT 2004

10/09823
PCT/JP03/04059

日 本 国 特 許 庁
JAPAN PATENT OFFICE

31.03.03

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出 願 年 月 日
Date of Application:

2002年 4月 4日

出 願 番 号
Application Number:

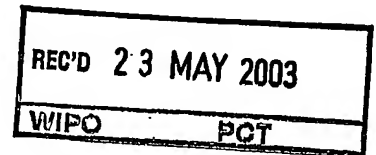
特願2002-102743

[ST.10/C]:

[JP2002-102743]

出 願 人
Applicant(s):

石原産業株式会社



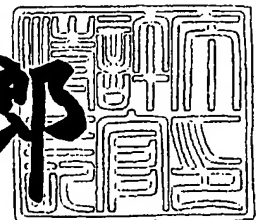
PRIORITY
DOCUMENT

SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)

2003年 5月 9日

特 許 庁 長 官
Commissioner,
Japan Patent Office

太田信一郎



出証番号 出証特2003-3033835

【書類名】 特許願

【整理番号】 181940

【提出日】 平成14年 4月 4日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/00

【発明者】

 【住所又は居所】 滋賀県草津市西湊川二丁目3番1号 石原産業株式会社
 中央研究所内

 【氏名】 石川 俊夫

【特許出願人】

 【識別番号】 000000354

 【住所又は居所】 大阪府大阪市西区江戸堀一丁目3番15号

 【氏名又は名称】 石原産業株式会社

【代理人】

 【識別番号】 100062144

 【弁理士】

 【氏名又は名称】 青山 蓀

【選任した代理人】

 【識別番号】 100086405

 【弁理士】

 【氏名又は名称】 河宮 治

【選任した代理人】

 【識別番号】 100098280

 【弁理士】

 【氏名又は名称】 石野 正弘

【手数料の表示】

 【予納台帳番号】 013262

 【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 データ解析装置および方法

【特許請求の範囲】

【請求項 1】 生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するデータ解析装置であって、

生体の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力手段と、

(1)説明変数を選択する選択手段と、(2)部分最小自乗法を実行して交差検証成績を計算する計算手段と、(3)前記(2)の計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定手段とを有し、(4)前記(1)の選択手段と前記(2)の計算手段と前記(3)の評価判定手段とを実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する決定手段とからなることを特徴とするデータ解析装置。

【請求項 2】 前記の選択手段において、説明変数を逐次取捨選択することを特徴とする請求項 1 に記載のデータ解析装置。

【請求項 3】 前記の選択手段において、遺伝的アルゴリズムを用いて説明変数を選択することを特徴とする請求項 1 に記載のデータ解析装置。

【請求項 4】 前記の計算手段において、1 個のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 1、2 または 3 に記載のデータ解析装置。

【請求項 5】 前記の計算手段において、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 1、2 または 3 に記載のデータ解析装置。

【請求項 6】 前記評価判定手段において、前記計算手段の結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと評価し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返すことを特徴とする請求項 1、2、4 または 5 に記載されたデータ解析装置。

【請求項 7】 前記の決定手段において、前記（１）の選択手段と前記（２）の計算手段と前記（３）の評価判定手段とを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定することを特徴とする請求項 1、2、3、4、5 または 6 に記載のデータ解析装置。

【請求項 8】 前記（１）の選択手段と前記（２）の計算手段とを複数のコンピュータで実行させることを特徴とする請求項 1、2、3、4 または 5 に記載のデータ解析装置。

【請求項 9】 請求項 1 で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力手段と、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定手段からなることを特徴とするデータ解析装置。

【請求項 10】 生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するデータ解析方法であって、

生体の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力ステップと、

（１）説明変数を選択する選択ステップと、（２）部分最小自乗法を実行して交差検証成績を計算する計算ステップと、（３）前記（２）の計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定ステップとを有し、（４）前記（１）の選択ステップと前記（２）の計算ステップと前記（３）の評価判定ステップとを実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する決定ステップとからなることを特徴とするデータ解析方法。

【請求項 11】 前記の選択ステップにおいて、説明変数を逐次取捨選択することを特徴とする請求項 10 に記載のデータ解析方法。

【請求項 12】 前記の選択ステップにおいて、遺伝的アルゴリズムを用いて説明変数を選択することを特徴とする請求項 10 に記載のデータ解析方法。

【請求項 13】 前記の計算ステップにおいて、１個のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 10、11 または 12 に記載のデータ解析方法。

【請求項 14】 前記の計算ステップにおいて、複数のサンプルを逐次除外

して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 10、11 または 12 に記載のデータ解析方法。

【請求項 15】 前記評価判定ステップにおいて、前記計算手段の結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと評価し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返すことを特徴とする請求項 10、11、13 または 14 に記載されたデータ解析方法。

【請求項 16】 前記の決定ステップにおいて、前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定することを特徴とする請求項 10、11、12、13、14 または 15 に記載のデータ解析方法。

【請求項 17】 前記(1)の選択ステップと前記(2)の計算ステップとを複数のコンピュータで実行させることを特徴とする請求項 10、11、12、13 または 14 に記載のデータ解析方法。

【請求項 18】 請求項 10 で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなることを特徴とするデータ解析方法。

【請求項 19】 生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定する、コンピュータにより実行されるデータ解析プログラムであって、

生体の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力ステップと、

(1)説明変数を選択する選択ステップと、(2)部分最小自乗法を実行して交差検証成績を計算する計算ステップと、(3)前記(2)の計算ステップの結果を評価し、説明変数の採用、不採用を判定する評価判定ステップとを有し、(4)前記(

1) の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する決定ステップとからなることを特徴とするデータ解析プログラム。

【請求項 2 0】 前記の選択ステップにおいて、説明変数を逐次取捨選択することを特徴とする請求項 1 9 に記載のデータ解析プログラム。

【請求項 2 1】 前記の選択ステップにおいて、遺伝的アルゴリズムを用いて説明変数を選択することを特徴とする請求項 1 9 に記載のデータ解析プログラム。

【請求項 2 2】 前記の計算ステップにおいて、1 個のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 1 9、2 0 または 2 1 に記載のデータ解析プログラム。

【請求項 2 3】 前記の計算ステップにおいて、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算することを特徴とする請求項 1 9、2 0 または 2 1 に記載のデータ解析プログラム。

【請求項 2 4】 前記評価判定ステップにおいて、前記計算手段の結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと評価し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返すことを特徴とする請求項 1 9、2 0、2 2 または 2 3 に記載されたデータ解析プログラム。

【請求項 2 5】 前記の決定ステップにおいて、前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定することを特徴とする請求項 1 9、2 0、2 1、2 2、2 3 または 2 4 に記載のデータ解析プログラム。

【請求項 2 6】 前記(1)の選択ステップと前記(2)の計算ステップとを複数のコンピュータで実行させることを特徴とする請求項 1 9、2 0、2 1、2

2または23に記載のデータ解析プログラム。

【請求項27】 請求項19で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなることを特徴とするデータ解析プログラム。

【請求項28】 上記の説明変数の選択において、初期状態では説明変数を全く含まないことを特徴とする請求項20に記載されたプログラム。

【請求項29】 上記の説明変数の選択において、初期状態では全説明変数を含むことを特徴とする請求項20に記載されたプログラム。

【請求項30】 上記の生体の状態が病気のタイプまたは医療診断の結果であることを特徴とする請求項19から29のいずれかに記載されたプログラム。

【請求項31】 請求項19から請求項30のいずれかに記載されたプログラムを記録した、コンピュータにより読み取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、生体の状態と遺伝子発現の量および／または細胞内物質の量との多変量解析処理に関するものである。

【0002】

【従来の技術】

2000年6月のヒトゲノムの解読宣言以降、ゲノムに書かれた遺伝情報がどのように発現して機能しているかを解明するポストゲノム時代に突入したと言われている。ヒトゲノム計画の進展の中で、ゲノム発現状態を測定する方法論も進展してきた。トランスクリプトーム(mRNA)測定手段としてオリゴヌクレオチドアレイやマイクロチップが知られている。またプロテオーム(蛋白質)測定手段として、以前からある2次元電気泳動に加えて、最近では質量分析の方法が進歩してきた。また抗体チップなどの先進の技術も注目されている。これらの測定技術は、生体の状態パラメータを短時間に一挙に測定できることがそれまでの技術と比較して画期的であるといえる。

【 0 0 0 3 】

遺伝子発現状態を効率的に測定する技術として次のものがあげられる。トランスクリプトーム (mRNA の総体) を特定するものとして、基盤に複数種の DNA を担持し、それに相補的な mRNA を検出する DNA チップが知られている。代表的な DNA チップには、遺伝子チップや DNA マイクロアレイがある。また、プロテオーム (蛋白質の総体) を特定するものには、2 次元電気泳動、抗体チップ、質量スペクトルを用いるものがある。またメタボローム (代謝中間体を含めた代謝産物の総体) を測定する手法も質量分析などによって試みられており、進展が見られる。

【 0 0 0 4 】

生体内の細胞の状態は遺伝子産物の発現によってよく記述されるため、従来の診断マーカーでは情報が不足している場面でも、精度のより高い診断が可能になるという期待も出てきている。たとえば、次のような研究があげられる。

【 0 0 0 5 】

P. O. Brown らは、DNA チップによってリンパ腫患者の細胞のトランスクリプトームを測定し、クラスター解析によって悪性と良性のリンパ腫 (DLBCL) を別クラスターに分離した (Nature 403(3), 503-11 (2000))。しかし、これは因果関係 (相関関係) のモデルを得る方法ではなく、どの遺伝子がどの程度重要かを判断できない。

【 0 0 0 6 】

A. Alaiya らは、2 次元電気泳動によって子宮がん患者 40 人の細胞のプロテオームを測定し、うち 22 人のデータから部分最小自乗法診断モデルを構築し、悪性度を説明した (Int. J. Cancer, 86, 731-36 (2000); Electrophoresis, 21, 1210-17 (2000); 国際公開 WO 00/70340)。その際、全変数モデルにおいて 1553 変数から loading の大きな 170 変数に限定することによって交差検証成績がよくなり ($Q^2 = 0.84$)、残り 18 患者の深刻度 (3 段階) を 11 / 18 の比率で正答した。交差検証法がモデル構築の際の指標になるという考えが表明されている。しかし、この方法では、loading を得る際にまず全変数モデルが成立しなければならない。また、それ以外の変数選択手法が考案されていない。

【0007】

J. Khanらは、DNAチップによって小児がん患者の細胞を測定し、ニューラルネットワークによって悪性度を説明した (Nature Medicine, 7(6), 673-79 (2001))。小児がん (SRBCT) 患者 88 人のトランスクリプトーム (6567 遺伝子) を測定し、うち 63 人のデータから主成分分析によって 10 次元に圧縮し、次に、人工ニューラルネットワーク診断モデルを構築した。ここで、影響力のある上位遺伝子を交差検証法によって絞り込み、96 遺伝子で最良の成績 (100%) を得た。このモデルで残り 25 人を予測し、93~100% の結果を得た。しかし、この方法でも、影響力を得る際にまず全変数モデルが成立しなければならない。またそれ以外の変数選択手法が考案されていない。10次元のような少ない変数の場合を扱えるが、変数の数が膨大な場合には適用できない。

【0008】

【発明が解決しようとする課題】

従来の診断マーカーでは情報が不足している場面でも、遺伝子発現情報を活用することで、より精度 (解像度) の高い診断が可能になるという期待も出てきている。遺伝子発現状態の測定結果は、膨大な情報量が得られることが従来にはなかった特徴であり、逆に情報量が多いために、効果的なデータ処理なくしてデータの活用はありえない。したがって、有用な知識を獲得するためには効果的な情報処理が欠かせない。前に説明したように、現状ではクラスター解析を中心とする方法が用いられているが、主成分分析などの方法も採用されている。クラスター解析や主成分分析は、教師付学習方法ではないため、病状の因果関係 (相関関係) のモデルを得ることはできない。すなわち、どの遺伝子がどの程度重要かを解析結果から得ることができないのが難点である。一方、部分最小自乗法は次元圧縮とモデルフィットを同時に行なう強力な多変量解析手法であるが、変数の数が膨大になった場合にしばしば有意な結果が得られない事態に直面する。したがって、膨大な遺伝子発現情報などから有用な知識を獲得できるような効果的な情報処理が望まれている。

【0009】

この発明の目的は、多変量の遺伝子発現情報、細胞内物質情報の効果的な情報

処理を提供することである。

【 0 0 1 0 】

【課題を解決するための手段】

本発明に係るデータ解析装置は、生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するデータ解析装置であって、生体の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力手段と、(1)説明変数を選択する選択手段と、(2)部分最小自乗法を実行して交差検証成績を計算する計算手段と、(3)前記(2)の計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定手段とを有し、(4)前記(1)の選択手段と前記(2)の計算手段と前記(3)の評価判定手段とを実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する決定手段とからなる。前記選択手段は、たとえば、説明変数を逐次取捨選択したり、遺伝的アルゴリズムを用いて説明変数を選択する。計算手段は、たとえば、1個のサンプルを逐次除外したり、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算する。評価判定手段は、たとえば、計算手段の結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと判定し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返す。決定手段は、たとえば、選択手段と計算手段と評価判定手段とを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する。また、選択手段と計算手段とを複数のコンピュータで実行させることもできる。こうして、相関モデルを構成するとき、多変量解析において、交差検証成績を基準に最適化させることにより説明変数を取捨選択し、説明変数の数を減らす次元圧縮をして良好なモデルを得る。

【 0 0 1 1 】

本発明に係るデータ解析方法は、生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するデータ解析方法であって、生体

の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力ステップと、(1)説明変数を選択する選択ステップと、(2)部分最小自乗法を実行して交差検証成績を計算する計算ステップと、(3)前記(2)の計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定ステップとを有し、(4)前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する決定ステップとからなる。前記選択ステップは、たとえば、説明変数を逐次取捨選択したり、遺伝的アルゴリズムを用いて説明変数を選択する。計算ステップは、たとえば、1個のサンプルを逐次除外したり、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算する。評価判定ステップは、たとえば、計算ステップの結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと判定し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返す。決定ステップは、たとえば、選択ステップと計算ステップと評価判定ステップとを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する。また、選択ステップと計算ステップとを複数のコンピュータで実行させることもできる。

【 0 0 1 2 】

本発明に係るデータ解析プログラムは、生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定する、コンピュータにより実行されるデータ解析プログラムであって、生体の状態を目的変数とし、複数の遺伝子発現の量および／または細胞内物質の量を説明変数とするサンプルの集合を入力する入力ステップと、(1)説明変数を選択する選択ステップと、(2)部分最小自乗法を実行して交差検証成績を計算する計算ステップと、(3)前記(2)の計算ステップの結果を評価し、説明変数の採用、不採用を判定する評価判定ステップとを有し、(4)前記(1)の選択ステップと前記(2)の計算ステップと前記(3)の評価判定ステップとを実行して部分最小自乗法モデルの交差検証成績を改

善し続けて部分最小自乗法モデルを決定する決定ステップとからなる。

【 0 0 1 3 】

上記のデータ解析プログラムにおいて、前記選択ステップは、たとえば、説明変数を逐次取捨選択したり、遺伝的アルゴリズムを用いて説明変数を選択する。計算ステップは、たとえば、1個のサンプルを逐次除外したり、複数のサンプルを逐次除外して部分最小自乗法を実行して交差検証成績を計算する。評価判定ステップは、たとえば、計算ステップの結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと判定し、説明変数を取捨選択しながら交差検証成績の評価判定を繰り返す。決定ステップは、たとえば、選択ステップと計算ステップと評価判定ステップとを繰り返し実行して部分最小自乗法モデルの交差検証成績を改善し続けて部分最小自乗法モデルを決定する。また、選択ステップと計算ステップとを複数のコンピュータで実行させることもできる。さらには、上記の説明変数の選択において、たとえば、初期状態では説明変数を全く含まないか、或いは、初期状態では全説明変数を含むこともできる。

【 0 0 1 4 】

上記のデータ解析プログラムにおいて、上記の生体の状態は、たとえば病気のタイプまたは医療診断の結果である。

【 0 0 1 5 】

また、本発明は、前記で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力手段と、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定手段からなるデータ解析装置、前記で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなるデータ解析方法及び前記で決定された相関モデル及び予測対象のサンプルについて当該モデルにおいて採用された説明変数を入力する入力ステップと

、入力された該説明変数に基づいて該サンプルの生体の状態を予測判定する予測判定ステップからなるデータ解析プログラムも包含する。

【0016】

本発明に係るコンピュータにより読取可能な記録媒体は、上記のいずれかのプログラムを記録する。

【0017】

【発明の実施の形態】

以下、添付の図面を参照して本発明の実施の形態を説明する。

以下に、選択された生体の状態と遺伝子発現の量および／または細胞内物質の量との相関モデルの決定について説明する。ここで、遺伝子発現の用語は、mRNA発現(トランスクリプトーム)や、mRNAによる翻訳の結果として生じる蛋白質(プロテオーム)を含むものとして用いる。また、細胞内物質の量とはここではたとえば、代謝中間体を含めた代謝産物全部であるメタボロームを意味する。たとえば、トランスクリプトーム(mRNA)やプロテオーム(蛋白質)の解析において、各サンプルデータは、生体の状態と遺伝子発現の量などからなる。各サンプルはたとえば1000個以上の膨大な遺伝子発現の量を含む。生体の状態は、たとえば病気のタイプまたは病気の診断指標であるが、より一般的には生体情報であればよい。「病気の診断指標」には、病気の進行度合いのほか、重篤度、深刻度などの表現で表わされるものも含む。ここで、遺伝子発現の量などの測定データは膨大な情報量からなるので、コンピュータを用いた効率的な多変量解析が必要である。

【0018】

データ収集において、予めいくつかのサンプルについて生体の状態(たとえば診断指標)を判定し、また、そのサンプルされたものから細胞液を獲得し、その細胞液中の多くの遺伝子産物の発現の量などを測定する。本発明の実施の形態のデータ解析では、こうして得られた遺伝子産物の発現の量などと生体の状態(たとえば診断指標)を入力し、相関モデル(たとえば部分最小自乗法モデル)を得る。ここで、コンピュータによる多変量解析プログラムを用いて、診断指標を目的変数とし、遺伝子発現の量および／または細胞内物質の量を説明変数とする因

果関係型の解析を行なって、各説明変数の重要性や影響度に関する情報を得る。

【0019】

我々は今回、遺伝子発現による医療診断という分野において、データ解析において、交差検証 (cross validation) の成績を最適化するように変数を選択することによって良好な相関モデル (たとえば部分最小自乗法モデル) が得られることを見出した。交差検証法では、手持ちのデータを複数群に分割し、その一部のデータ群 (訓練集合) だけを使ってフィットしたモデルを用いて残る別のデータ群 (テスト集合) を予測することによって、モデルの予測力を試す。通常の部分最小自乗法 (PLS) においては潜在変数の次元選択に交差検証法が用いられているが、ここでは、部分最小自乗法において、潜在変数を1次元に固定し、1以上の入力変数 (説明変数) を逐次取捨選択しながら、交差検証成績 (たとえば平方和の予測誤差) を最適化した。その結果、全変数を採用した場合には有意な相関モデルを得られなかった場合にも、良好でかつ予測力のある相関モデルが得られることが判明したのである。この交差検証法を用いた変数選択の逐次取捨選択により、安定な相関モデルが得られる。なお、ここでいう「最適化」とは、交差検証成績が、説明変数を取捨選択するための、そのときの解析条件の範囲で、改善がみられなくなるまで改良したことを意味しており、交差検証成績がすべての説明変数の組合せの中で最適なものを見出したという意味ではない。この変数選択手法を用いると、病状を決定する因子を少数に特定し、廉価な診断用材料 (DNAチップ、抗体チップ、DNA含有ベクターなど) を設計でき、それ自体独自の価値を持つものである。また、この変数選択手法は、予め設定される各種の変数選択条件と共に運用することが可能である。

【0020】

上に述べたように、説明変数は、交差検証成績を基準に逐次取捨選択される。説明変数を追加する場合は、その説明変数について、交差検証成績評価が改善されなかったと判定された場合には当該説明変数を除外し、改善されたと判定された場合には当該説明変数を追加する。また、説明変数を除外する場合は、その説明変数について、交差検証成績評価が改善されなかったと判定された場合には当該説明変数を除外せず、改善されたと判定された場合には当該説明変数を除外す

る。ここで、1以上の説明変数を選択した場合に、交差検証成績評価は次のように進める。n個のサンプルからいくつかのサンプルを逐次除外して部分最小自乗法モデルを求め、各モデルにおいて除外したサンプルの遺伝子発現の量から予測される生体の状態を示す目的変数と、除外したサンプルの生体の状態を示す目的変数との各々の誤差の代表値を求める。「代表値」とは、和、平均、最大値、中位値、最頻値などのデータを特徴づける値をいう。そして、当該誤差の代表値が小さくなった場合に、交差検証成績が改善されたと判定し、当該説明変数を追加または削除する。この交差検証成績評価を、説明変数を取捨選択しながら逐次繰り返して、交差検証成績を改善し続ける。改善されなくなれば交差検証成績を最適化したとして説明変数の取捨選択を終了する。その結果、取捨選択により絞り込んだ数の説明変数からなる最適な部分最小自乗法モデルが得られる。

【0021】

因果関係型の解析手法においてはオーバーフィット (over fitting) を避けるための工夫が必要となる。ここでいうオーバーフィットとは、説明変数が多すぎるためにたまたま予測結果と実績とが一致するものの、本当の相関関係をとらえ損なっているため、モデルフィットに用いたデータ以外に予測能力を持たないことをいう。ここでは、相関モデルとして部分最小自乗法を用いるが、部分最小自乗法は次元圧縮とモデルフィットを同時に行なう強力な多変量解析手法であり、オーバーフィットの問題に比較的強いとされている。しかし遺伝子発現状態解析のように膨大な変数を扱う場合には、有意な結果が得られない事態に直面する。従来技術として説明したAlaiyaやKhanの手法は全変数モデルが有意に成立することを前提としているので、変数の絞り込みには一般的には適用できない。これに対し、本発明では、交差検証予測結果を最適にするように変数を絞り込むことにより、オーバーフィットを減らすことができた。また、本発明は、前記Khanの手法とは異なり、主成分分析などの前処理を介さない方法である。従来技術では、説明変数が膨大な場合には、有意なモデルを得ることができないことから、予め、全説明変数を基にたとえば、主成分分析などで次元圧縮する前処理をし、これによって得られた説明変数によって解析する方法が用いられる。しかしながら、この方法では、構成したモデルで予測を行なうためには、モデル構成の基となった

全説明変数が必ず必要となり、たとえば、説明変数が遺伝子発現の量であれば、診断用遺伝子チップに担持する遺伝子としては、モデル構成に用いた遺伝子の全てが必要となるかもしくは別の手法を用いて変数選択することが必要となる。一方、本発明においては、説明変数の選択によって説明変数を絞り込んでいるので、たとえば、説明変数が遺伝子発現の量であれば、診断用遺伝子チップに担持する遺伝子は、選択された説明変数に相当する遺伝子を担持すれば良いことになる。

【 0 0 2 2 】

なお、Todeschiniらは、有機化合物の大気中の分解を予測するため、遺伝的アルゴリズムによって交差検証成績を最適化するように変数選択を行ない、重回帰モデルを得ている (P. Gramatica, V. Consonni & R. Todeschini, Chemosphere 38(5), 1371-78 (1999))。53化合物と175記述子でモデル構築を行ない ($Q^2 = 0.79$)、7変数が選択され、98化合物の予測を行なった ($Q^2 = 0.75$)。交差検証成績を最適化するように変数選択を行なっている点では、本実施形態と同様の手法である。しかし、重回帰モデルを採用しているために、結果として選択される変数は少数個にとどまり、複数の遺伝子発現の量および／または細胞内物質の量の解析には適用できない。

【 0 0 2 3 】

本実施形態では、生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するとき、交差検証成績を基準に最適化させるように説明変数を逐次追加・除外することによって、説明変数の数を減らし、良好な相関モデルを得る。このようなアプローチの優位性は、下記の実施例から推測されるように、次のとおりである。

- 1) 病気や生体现象の背後で働いている重要な遺伝子やメカニズムを推定／特定でき、理解が深まる。
- 2) 重要な遺伝子産物や細胞内物質だけに絞った廉価な診断用材料 (DNAチップ、抗体チップなど) の設計が可能になる。

【 0 0 2 4 】

本実施形態では、交差検証成績を最適化するように説明変数を段階的に取捨選

択するが、たとえば具体的には、ステップワイズ(step wise)法に代表される説明変数を選択する選択手段と、リーブ・ワン・アウト(leave-one-out)法に代表される交差検証法に部分最小自乗法を適用して計算する計算手段と、前記計算手段の結果を評価し、説明変数の採用、不採用を判定する評価判定手段とを組合せて用いる。すなわち、 m 個の説明変数の中から1以上の説明変数を選択し、次いで、部分最小自乗法を実行して交差検証成績を計算し、更に、該計算結果を評価して前記で選択した説明変数の採用、不採用を判定する。前記評価判定では、前記計算手段の結果から、各計算において除外したサンプルの遺伝子発現から予測される生体の状態を示す目的変数値と、前記除外したサンプルの生体の状態を示す目的変数値との誤差の代表値を求め、当該誤差の代表値が小さくなった場合に、その交差検証成績が改善されたと評価し、説明変数の取捨選択を判定する。上記したように、選択手段と計算手段と評価判定手段とを実行して部分最小自乗法モデルの交差検証成績を改善し続けて、その改善がみられなくなるまで改良し、部分最小自乗法モデルを決定する。なお、本実施形態では、サンプルを1個ずつ逐次除外している(リーブ・ワン・アウト法)が、その代わりに、複数のサンプルを除外して交差検証成績を評価してもよい(リーブ・ n ・アウト法)し、また、Khan et al.により用いられた3分割法(three-fold)等の他の方法を用いることもできる。3分割法では、説明変数をランダムにシャッフルして3つのグループに分ける。その中の2つのグループを用いてモデルを構成し、残りの1つのグループでモデルを評価する。また、説明変数の選択方法としてはステップワイズ法、非線形アルゴリズム(たとえば遺伝的アルゴリズムなど)を用いてもよく、変数選択に関して予め何らかの条件が分っていれば、それに応じて探索範囲を限定できる。

【0025】

次に、データの収集と解析について具体的に説明する。図1は、遺伝子発現解析システムを示す。データ収集のため、予めいくつかのサンプルについて診断指標(たとえば病気のタイプないし進行度合いを含む)を判定し、また、そのサンプルされたものから細胞液を獲得し、DNAチップを用いてその細胞液中の多くの遺伝子産物の発現の量を測定する。測定には、共焦点型レーザスキャナ(たと

例えばAflymetrix社、428アレイスキャナ) 10を用いる。吸光度によりmRNAの量が測定される。このデータ収集は公知の方法である。測定データは、コンピュータ12に送られ解析される。コンピュータ12は、CPU14を備えた通常の構成のコンピュータであり、それに接続される記憶装置(たとえばハードディスク装置) 16の記録媒体(たとえばハードディスク)には、測定データ18や解析ソフト20が格納される。この解析ソフト20を用いてデータ18が解析され、生体の状態と遺伝子発現の量などとの相関モデルが決定される。

【0026】

図2は、コンピュータ12により実行される、生体の状態と遺伝子発現の量などとの相関モデルを得るためのデータ解析ソフト20のフローチャートを示す。まず、相関モデル作成用のデータを入力する(S10)。データはたとえばDNAチップを用いて収集したものである。入力データ(サンプル集合)は、それぞれ目的変数(たとえば診断指標)とm個(たとえば2000個)の説明変数(たとえば遺伝子発現の量)からなる。また、場合によっては、上述のデータ(訓練集合)以外に、テスト集合のデータを入力する。ここでテスト集合とは交差検証の評価のためのデータ群を意味するのではなく、モデル決定が終了した後にモデルの予測力をテストするためのデータ群である。

【0027】

まず、初期設定として、選択された説明変数の数を0とし、交差検証成績CVの最良値 CV_0 を $-\infty$ とする(S12)。次に、説明変数の選択を行う。まず、説明変数を指す番号iを1とし(S14)、第i変数(遺伝子発現の量)を仮に採用して(S16)、部分最小自乗法を実行し、交差検証成績CVを計算する(S18、図3参照)。ここで、リーブ・ワン・アウト処理を用いる。これは、たとえば50個のサンプルからなる訓練集合において、1番から50番の全てを順次1個ずつ除いて残りの49個のサンプルで予測した結果と、その時除いた1個の結果とを比較し、その誤差が大きい場合に、仮に選択した説明変数(第i変数)が適していないと判断する手法である。もし、得られた成績CVが現在の最良値 CV_0 より最適化されれば(S20でYES)、第i変数を採用し、かつ、成績CVを新しい最良値 CV_0 に更新する(S22)。しかし、得られた成績CVが最良値 CV_0 より

大きくなければ (S 2 0 で NO)、第 i 変数を採用しない (S 2 4)。そして、ステップ S 1 4 に戻り、同様の処理を繰り返す。この処理を交差検証成績 CV が改善されなくなる (S 2 6 で NO) まで繰り返す。ここで、相関モデルに採用する説明変数については 1 つずつ段階的に増加 (追加) または減少 (除外) して成績 CV を評価判定している。すなわち、全体としての合致度合いがよくなるように各説明変数を解析に加えるかどうかを逐次判定しながら、説明変数の取捨選択を行い、これを、全体としての合致度合いがよくなるまで繰り返す。以上の処理で改善があると、ふたたびステップ S 1 4 の初め ($i=1$) に戻り、それまでに選択されている説明変数を基に、さらに説明変数の選択を繰り返す。なお、ここではモデルの予測力を判断するために、訓練集合とテスト集合とに予め分割しておいたデータ集合を用いてデータ解析しており、上述の解析は、訓練集合を用いて行なった結果であるので、この結果からテスト集合について予測を行い、実測データとの合致度を評価 (S 2 8) している。このような評価は必ずしも必要でないが、予測力を判断するには有効である。

【 0 0 2 8 】

図 3 は、リーブ・ワン・アウト処理を含む交差検証成績 CV の計算 (図 2、S 1 8) のフローチャートを示す。ここで、選択された変数について交差検証成績が計算される。まず、PRESS の初期値を 0 とする (S 1 8 0)。次に、 n 個の集合内のサンプルを指す番号 j を 1 とし (S 1 8 2)、第 j サンプル以外の $n-1$ 個のサンプルで部分最小自乗法を実行し (S 1 8 4)、第 j サンプルの目的変数を予測する (S 1 8 6)。差の自乗を計算して PRESS に加算する (S 1 9 0)。次に番号 j を 1 増加し (S 1 8 2)、同様の処理をおこなう。これを番号 $j=n$ まで各サンプルについて繰り返す。得られた PRESS は、1 個のサンプルを順次除外して計算した予測値と実測値との差の平方和であり、予測誤差を表わす量である。この予測残差自乗和 PRESS の符号を変えたものを交差検証成績 CV とする (S 1 9 2)。

【 0 0 2 9 】

本実施形態では、交差検証法を用いて、入力変数 (説明変数) を段階的に 1 つずつ追加・除外しながら、交差検証成績 (平方和の予測誤差) を最適化する。ここ

で、説明変数の段階的な追加・除外の内容を理解しやすくするため、以下で、さらに具体的に5つのモデル構築手法について説明する。これらは、説明変数の逐次的な選択の手順が異なる。

【0030】

図4は、第1のモデル構築手法を示す。データ集合においてどの説明変数も選択されていない状態を初期状態とする（S112）。次に、1番目の説明変数から最後(m番目)の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ（S118）を繰り返しながら判定し、改善する場合にはその説明変数を追加する（S114～S124）。そのような改善と追加がなくなる（S126でNO）まで、1番目の説明変数から上記逐次判定操作を繰り返す。

【0031】

さらに詳しく説明すると、まず、初期設定として、選択された説明変数の数 N_P を0とし、交差検証成績CVの最良値 CV_0 を $-\infty$ とする（S112）。次に、説明変数の選択を行う。まず、変数 i を1とし（S114）、第 i 変数を仮に採用する（S116）。ただし、第 i 変数がすでに採用されていれば（S115でYES）、ステップS114に戻る。次に、部分最小自乗法を実行し、交差検証成績CVを計算する（S118）。ここで、リーブ・ワン・アウト処理を用いる。もし、得られた成績CVが現在の最良値 CV_0 より最適化されれば（S120でYES）、第 i 変数を採用し、かつ、成績CVを新しい最良値 CV_0 に更新する（S122）。しかし、得られた成績CVが最良値 CV_0 より大きくなければ（S120でNO）、第 i 変数を採用しない（S124）。そして、ステップS114に戻り、同様の処理を繰り返す。この処理を交差検証成績CVが改善されなくなる（S126でNO）まで繰り返す。以上の処理で改善があると、ふたたびステップS114に戻り、新しいループを開始する。ここで、それまでに選択されている変数を基に、さらに変数の選択を繰り返す。こうして、データ集合を用いて選択された変数を用いた相関モデルが得られる。

【0032】

図5は、第2のモデル構築手法を示す。この手法では、全ての説明変数が選択されている状態を初期状態とする（S212）。次に、1番目の説明変数から最後（m番目）の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ（S218）を繰り返しながら判定し、改善する場合にはその説明変数を除外する（S214～S224）。そのような改善と除外がなくなる（S226でNO）まで、1番目の説明変数から上記逐次判定操作を繰り返す。

【0033】

さらに詳しく説明すると、まず、初期設定として、選択された説明変数の数NをPとし、交差検証成績CVの最良値 CV_0 を $-\infty$ とする（S212）。すなわち、すべての説明変数を選択する。次に、説明変数の選択を行う。まず、変数iを1とし（S214）、第i変数を仮に除外する（S216）。ただし、第i変数がすでに除外されていれば（S215でYES）、ステップS214に戻る。部分最小自乗法を実行し、交差検証成績CVを計算する（S218）。ここで、リーブ・ワン・アウト処理を用いる。もし、得られた成績CVが現在の最良値 CV_0 より最適化されれば（S220でYES）、第i変数を除外し、かつ、成績CVを新しい最良値 CV_0 に更新する（S222）。しかし、得られた成績CVが最良値 CV_0 より大きくなければ（S220でNO）、第i変数を除外しない（S224）。そして、ステップS214に戻り、同様の処理を繰り返す。この処理を交差検証成績CVが改善されなくなる（S226でNO）まで繰り返す。以上の処理で改善があると、ふたたびステップS214に戻り、新しいループを開始する。ここで、それまでに選択されている変数を基に、さらに変数の選択を繰り返す。こうして、データ集合を用いて選択された変数を用いた相関モデルが得られる。

【0034】

図6は、第3のモデル構成手法を示す。この手法は、第1と第2の手法の直列的な組合せである。まず、どの説明変数も選択されていない状態を初期状態とする（S112）。次に、1番目の説明変数から最後（m番目）の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差

検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を追加選択し、そのような改善と追加がなくなるまで1番目の説明変数から上記逐次判定操作を繰り返す（S114～S126）。次に、1番目の説明変数から最後（m番目）の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を除外し、そのような改善と除外がなくなるまで、1番目の説明変数から上記逐次判定操作を繰り返す（S214～S226）。

【0035】

図7は、第4のモデル構築手法を示す。この手法は、第3の手法の変形である。まず、どの説明変数も選択されていない状態を初期状態とする（S112）。次に、1番目の説明変数から最後（m番目）の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ（S118）を繰り返しながら判定し、改善する場合にはその説明変数を追加選択する（S114～S124）。そのような改善と追加がなくなる（S126でNO）まで、1番目の説明変数から上記逐次判定操作を繰り返す。次に、1番目の説明変数から最後（m番目）の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ（S218）を繰り返しながら判定し、改善する場合にはその説明変数を除外する（S214～S224）。そのような改善と除外がなくなる（S226でNO）まで、1番目の説明変数から上記逐次判定操作を繰り返す。上記逐次判定追加改善ステップまたは上記逐次判定除外改善ステップで少なくとも一度改善があれば（S227でYES）、ステップS112に戻り、上記操作（S112～S227）を繰り返す。これを改善がなくなる（S227でNO）までおこなう。

【0036】

図8は、第5のモデル構築手法を示す。この手法は、第1と第2のスキームの

並列的な組合せである。どの説明変数も選択されていない状態を初期状態とする (S112)。次に、1番目の説明変数から最後 (m番目) の説明変数までの説明変数ごとに逐次、その説明変数が選択されていない場合にはその説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ (S118) を繰り返しながら判定し、改善する場合にはその説明変数を追加する (S114~S124)。また、選択する説明変数ごとに、その説明変数がすでに選択されている場合には、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理を用いた交差検証成績評価ステップ (S218) を繰り返しながら判定し、改善する場合にはその説明変数を除外する (S216~S224)。そのような改善と追加または除外がなくなる (S126でNO) まで、1番目の説明変数から上記逐次判定操作を繰り返す。

【0037】

次に、第4のモデル構築手法(図7)を適用した場合を、表1のデータ集合を例として説明する。このデータ集合に対して、部分最小自乗法による解析を用いて相関モデルを求める。表1のデータでは、サンプルの数 n は10であり、また、説明を容易にするため、説明変数の数 m は19と少なくしている。表1において、 p_1 は目的変数を表わし、 $p_2 \sim p_{20}$ は説明変数を表わす。(ただし表1では、表示の便宜のため、 p_{16} 以降のデータを省略している。)第4手法(図7)のステップ S114、S214とは異なり、説明変数を表わす i は p_{20} から p_2 まで逆に逐次処理することとした。CV評価値としてここでは予測残差自乗和(PRESS)を採用した。PRESSが小さいほど、CV評価値はよい。初期状態では、採用された説明変数の数 NP は0であり、 $PRESS = \infty$ ($CV_0 = -\infty$) である。

【0038】

【表 1】

表 1 10個のサンプルのデータ

#	p1	p2	p3	p4	p5	p6	p7	p8	p9	p10	p11	p12	p13	p14	p15
1	0.713	0.105	0.782	0.425	0.164	0.023	0.696	0.543	0.333	0.691	0.336	0.668	0.017	0.061	0.5
2	0.133	0.009	0.071	0.002	0.793	0.872	0.092	0.391	0.63	0.241	0.517	0.369	0.166	0.841	0.1
3	0.545	0.193	0.765	0.334	0.109	0.538	0.578	0.652	0.38	0.501	0.729	0.91	0.865	0.389	0.8
4	0.752	0.915	0.472	0.999	0.798	0.363	0.622	0.487	0.353	0.967	0.778	0.484	0.517	0.982	0.0
5	0.9	0.407	0.534	0.816	0.806	0.42	0.572	0.957	0.12	0.696	0.833	0.051	0.377	0.849	0.4
6	0.455	0.587	0.721	0.53	0.252	0.434	0.882	0.486	0.741	0.243	0.893	0.947	0.462	0.952	0.2
7	0.427	0.652	0.515	0.426	0.764	0.592	0.595	0.595	0.551	0.606	0.416	0.163	0.316	0.718	0.6
8	0.042	0.902	0.274	0.899	0.402	0.469	0.668	0.945	0.746	0.912	0.97	0.515	0.368	0.514	0.4
9	0.935	0.276	0.936	0.101	0.54	0.356	0.899	0.71	0.924	0.792	0.486	0.329	0.501	0.076	0.5
10	0.54	0.021	0.505	0.224	0.724	0.431	0.071	0.968	0.482	0.322	0.773	0.543	0.353	0.107	0.9

【0039】

【表 2】

表 2 表 1 のデータについての 10 の段階での変数選択結果

0		∞	-
1	追加	p20	0.111 p20
2	追加	p18	0.090 p18 & p20
3	追加	p16	0.073 p16 & p18 & p20
4	追加	p10	0.073 p10 & p16 & p18 & p20
5	追加	p6	0.062 p6 & p10 & p16 & p18 & p20
6	追加	p3	0.060 p3 & p6 & p10 & p16 & p18 & p20
7	追加	p12	0.055 p3 & p6 & p10 & p12 & p16 & p18 & p20
8	除外	p20	0.053 p3 & p6 & p10 & p12 & p16 &
9	除外	p10	0.050 p3 & p6 & p12 & p16 & p18
10	追加	p13	0.048 p3 & p6 & p12 & p13 & p16 & p15

【0040】

先に述べたように、変数はp20からp2まで逆の順で処理する。表2は、表1のサンプルについて、左端の数字は、変数の取捨選択で改善がみられた10の段階を示す。なお、0は初期状態を意味する。次の列の「追加」と「除外」は、追加のループと除外のループの処理であることを意味する。次の列の変数は、追加または除外された変数を示す。次の列は、交差検証成績(PRESSをサンプル数で割ったもの)を示す。右端の列は、その段階で選択されている変数を示す。

【0041】

初期状態では、変数は全くない状態であり、PRESSは ∞ である。表2に示すように、最初、p20を説明変数として採用すると、PRESS=0.111となり、初期値に比べて改善されるので、説明変数p20の追加を実施する。次に、変数p19を加えてp19とp20の2つを説明変数とすると、PRESS=0.129となり改善をもたらさないので、p19は追加しない。次に、説明変数p18を加えるとPRESS=0.090となり、改善するので、p18を追加し、p18とp20を説明変数とする。以下同様に表2に示すように続く。(ここで、p10を追加採用するのは、小数点以下4桁目で改善されているためである。)説明変数p20～p2の1回目のループを終了した時点で、説明変数がp3、p6、p10、p16、p18およびp20となり、PRESS=0.60となる。2回目のループでは、説明変数p12が追加され、PRESS=0.55となる。3回目のループでは追加による改善がなく、ひとまずS114～S126の追加処理を終了し、S214に移る。この時点での部分最小自乗法のフィットならびにリーブ・ワン・アウト予測状況は表3のとおりである。

【0042】

表3は、10のサンプルについて、表2の7で示す段階まで処理が進んだ時点での部分最小自乗法のフィットならびにリーブ・ワン・アウト予測状況を示す。ここで、モデル予測とリーブ・ワン・アウト予測のそれぞれにおいて、計算値と実測値との誤差を示す。さらに、その下側に、誤差の自乗平均、相関係数Rの自乗および予測相関係数Qの自乗を示す。

【0043】

【表 3】

表 3 表 2 の段階 7 での処理結果

#	モデル予測値			リーフ・ランアウト予測	
	実測値	計算値	誤差	計算値	誤差
1	0.713	0.757	-0.044	0.693	0.020
2	0.133	-0.056	0.189	-0.051	0.184
3	0.545	0.497	0.048	0.480	0.065
4	0.752	0.646	0.106	0.495	0.257
5	0.900	0.687	0.214	0.557	0.343
6	0.455	0.489	-0.034	0.512	-0.057
7	0.427	0.624	-0.198	0.672	-0.245
8	0.042	0.349	-0.307	0.517	-0.475
9	0.935	0.865	0.070	0.782	0.153
10	0.154	0.197	-0.044	0.285	-0.132
<hr/>					
	0.093		0.024		0.055
	$R^2=0.744$		$Q^2=0.407$		

【 0 0 4 4 】

次に、S 2 1 4 から始まる除外処理の 1 回目のループにおいて、説明変数 p10 と p20 を除外することが改善をもたらした。2 回目のループでは改善がなく、S 2 1 4 ～ S 2 2 6 を終了するが、S 2 2 7 の判断により再度 S 1 1 2 に戻る。次に、追加処理の 1 回目のループにおいて、p13 の追加だけが改善をもたらしたが、続く除外処理の 1 回目のループでは、改善がなかった。もう一度 S 1 1 2 に戻り、ステップ S 1 1 4 ～ S 1 2 6 およびステップ S 2 1 4 ～ S 2 2 6 では改善がなくなったのを確認して、処理を終了した。こうして選択された説明変数は、p3、p6、p12、p13、p16 および p18 の 5 個であり、PRESS=0.048 となった。詳細は表 4 のとおりである。

【 0 0 4 5 】

表4は、表2の段階10まで処理が進んだ時点での部分最小自乗法のフィットならびにリーブ・ワン・アウト予測状況を示す。

【0046】

【表4】

表4 表2の段階10での処理結果

#	実測値	モデル予測		リーブ・ワン・アウト予測	
		計算値	誤差	計算値	誤差
1	0.713	0.771	-0.058	0.663	0.050
2	0.133	-0.013	0.146	0.041	0.092
3	0.545	0.610	-0.065	0.595	-0.050
4	0.752	0.524	0.228	0.380	0.372
5	0.900	0.696	0.205	0.543	0.357
6	0.455	0.591	-0.137	0.623	-0.168
7	0.427	0.638	-0.211	0.696	-0.269
8	0.042	0.189	-0.147	0.268	-0.226
9	0.935	0.841	0.094	0.756	0.179
10	0.154	0.209	-0.055	0.294	-0.140
		0.093	0.022	0.048	
		$R^2=0.765$		$Q^2=0.482$	

【0047】

なお、説明変数の数が多い時に強いとされる部分最小自乗法であるが、p20～p2の全てを説明変数として採用した場合には、表5に示すように、PRESS=0.124となった。すなわち、リーブ・ワン・アウト処理は、平均値からの誤差(0.093)よりも悪い成績をもたらす。

【0048】

【表 5】

表 5 全ての説明変数を採用した場合の処理結果

#	モデル予測			リーブ・ワン・アウト予測	
	実測値	計算値	誤差	計算値	誤差
1	0.713	0.712	0.001	0.527	0.186
2	0.133	-0.073	0.206	0.222	-0.090
3	0.545	0.561	-0.016	0.538	0.007
4	0.752	0.656	0.096	0.351	0.402
5	0.900	0.691	0.209	0.432	0.469
6	0.455	0.519	-0.064	0.562	-0.107
7	0.427	0.583	-0.156	0.629	-0.203
8	0.042	0.430	-0.388	0.724	-0.682
9	0.935	0.794	0.140	0.480	0.454
10	0.154	0.182	-0.029	0.457	-0.303
<hr/>					
	0.093		0.029	0.124	
	$R^2=0.684$			$Q^2=-0.330$	

【0049】

次に、実施例について説明する。P.O. Brownらのホームページ (<http://llmp.p.nih.gov/lymphoma/>) より入手した28名のDLBCL（リンパ腫）患者のデータを、20名のデータからなる訓練集合と8名のデータからなるテスト集合に分けた。目的変数に生存月数を採用し、説明変数には18432スポットのうち、28データにおいてch1、ch2ともに正の数となる12832スポットの $\log(ch1/ch2)$ 値を採用した。

【0050】

訓練集合において部分最小自乗法（PLS）のモデル決定を試みた。12832変数全てを用いて部分最小自乗法の解析をしたところ、リーブ・ワン・アウト予測は有意($Q^2 > 0.5$)にはならなかった。次にリーブ・ワン・アウト予測誤

差が最小になるように説明変数を段階的に1つつ増減した。モデル構成手法としては前述の第3のモデル構成手法において説明変数の追加及び除外の順番並びにリーブ・ワン・アウト処理におけるサンプルの除外の順番が異なるほかは同様な方法を用いた。すなわち、どの説明変数も選択されていない状態を初期状態とする(S112)。次に、最後(m番目)の説明変数から最初(1番目)の説明変数までの未だ選択されていない説明変数ごとに逐次、その説明変数を選択した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理(ここでは、最後(n番目)のサンプルから最初(1番目)のサンプルを逐次除外した)を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を追加選択し、そのような改善と追加がなくなるまでm番目の説明変数から上記逐次判定操作を繰り返す(S114~S126)。次に、最後(m番目)の説明変数から最初(1番目)の説明変数までの選択されている説明変数ごとに逐次、その説明変数を除外した場合に交差検証成績が改善するかどうかを、リーブ・ワン・アウト処理(ここでも最後(n番目)のサンプルから逐次除外した)を用いた交差検証成績評価ステップを繰り返しながら判定し、改善する場合にはその説明変数を除外し、そのような改善と除外がなくなるまで、最後(m番目)の説明変数から上記逐次判定操作を繰り返す(S214~S226)。その結果、有意なモデル($R^2=0.988$ 、 $Q^2=0.895$ 、 $N_p=342$)を得た。図9は、このデータについての最小自乗法成績を示す。図9において、ひし形は訓練集合のデータ(20人)を示し、3角は、それらについての交差検証成績のデータを示す。また、4角はテスト集合のデータ(8人)を示す。得られた部分最小自乗法モデルは、テスト集合のうち、4/8をきわめて良好に、また1/8を良好に予測するものであった。

【0051】

なお、上述の多変量解析によるデータ解析では、扱ったサンプルはDNAチップを用いて得たデータであった。しかし、このデータ解析は、DNAチップを用いて得たデータに限定されるものではなく、蛋白質発現量、細胞内物質の量などのデータに対しても有用であろうことは容易に推測されることである。

【0052】

【発明の効果】

生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するとき、説明変数の選択と交差検証法とを用いて変数を絞り込むことができる。これにより、良好でかつ予測力のある多変量解析モデル（相関モデル）が得られる。特に遺伝子発現の量のように、説明変数の数がたとえば1000以上と膨大な場合に有用である。変数の数を少なくすることにより、病気や生体现象の背後で働いている重要な遺伝子やメカニズムを推定／特定でき、理解が深まる。また、重要な遺伝子産物や細胞内物質だけに絞った廉価な診断用材料（DNAチップ、DNA含有ベクター、抗体チップなど）を設計し、提供できる。

【図面の簡単な説明】

【図1】 遺伝子発現解析システムのブロック図

【図2】 解析ソフトのフローチャート

【図3】 交差検証成績CVの計算のフローチャート

【図4】 変数選択の第1モデル構築手法のフローチャート

【図5】 変数選択の第2モデル構築手法のフローチャート

【図6】 変数選択の第3モデル構築手法のフローチャート

【図7】 変数選択の第4モデル構築手法のフローチャート

【図8】 変数選択の第5モデル構築手法のフローチャート

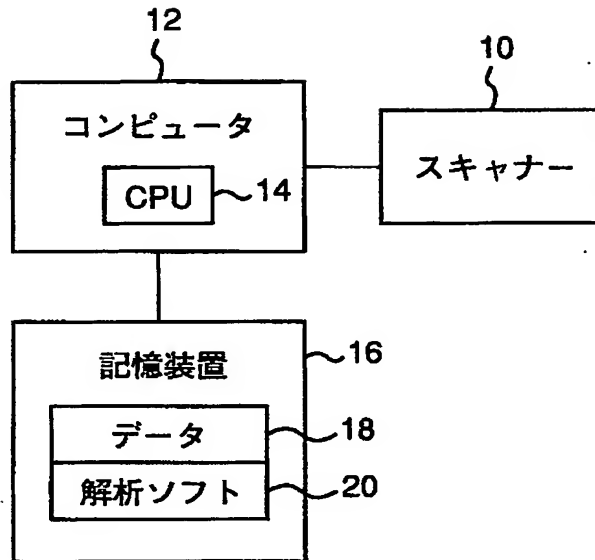
【図9】 最小自乗法モデルの成績を示すグラフ

【符号の説明】

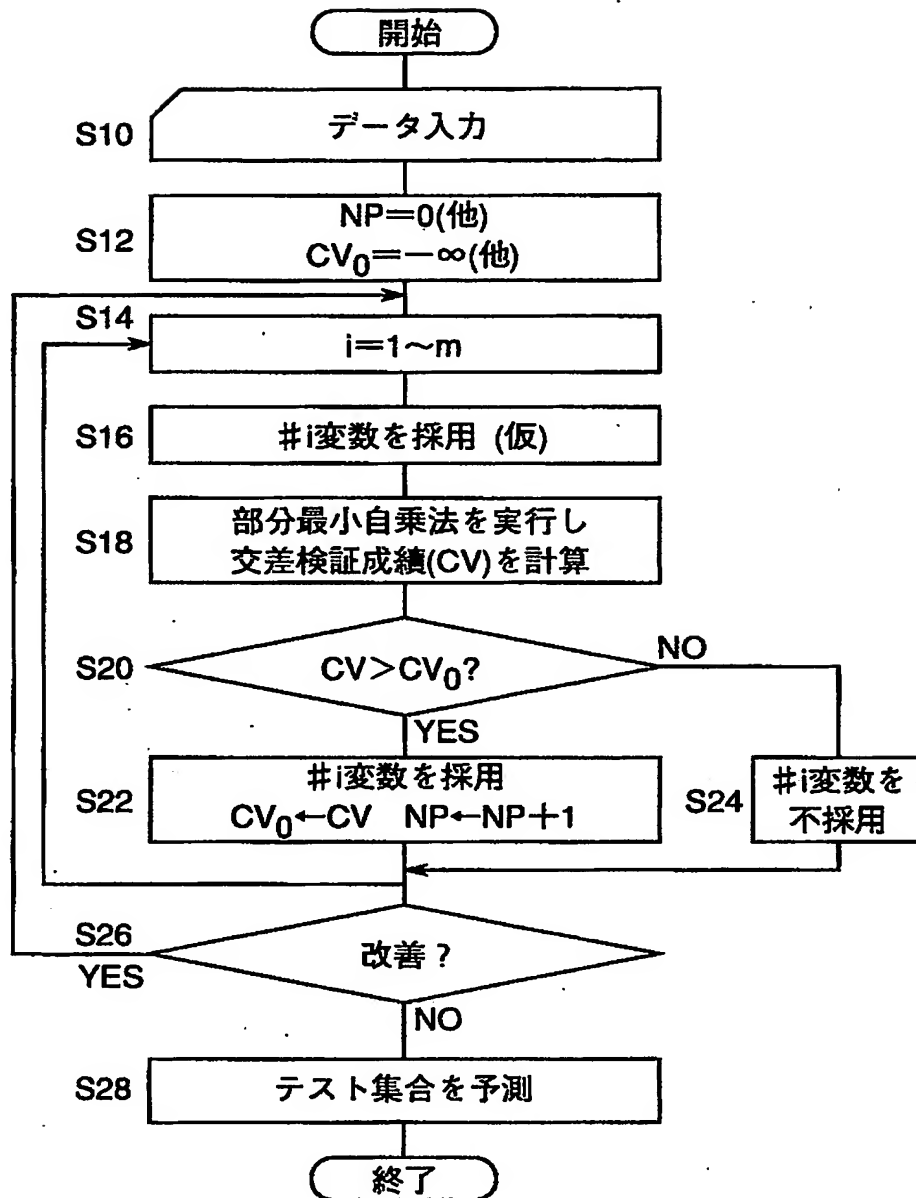
10 スキャナ、 12 コンピュータ、 14 CPU、 16 記憶装置、 18 記録媒体、 20 解析ソフト。

【書類名】 図面

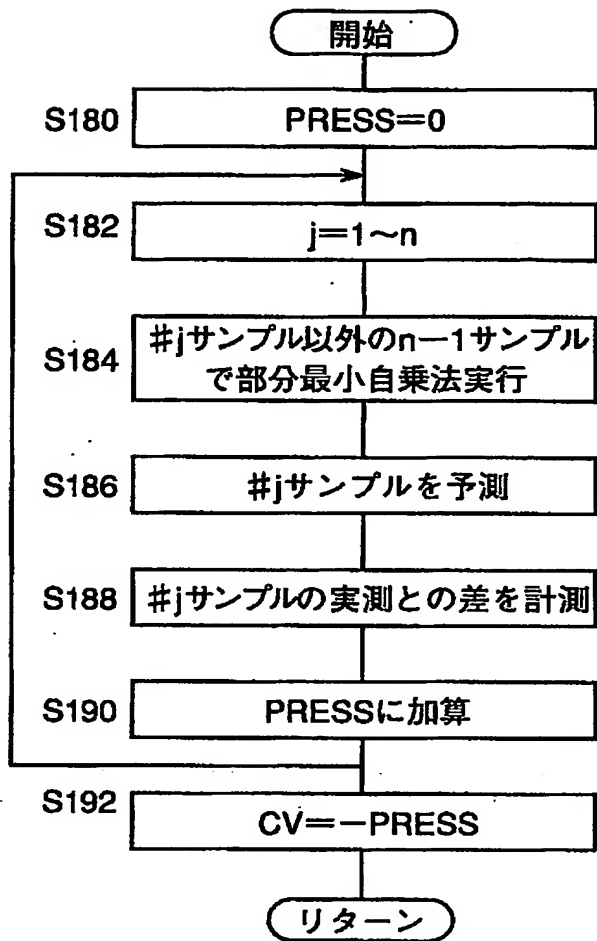
【図 1】



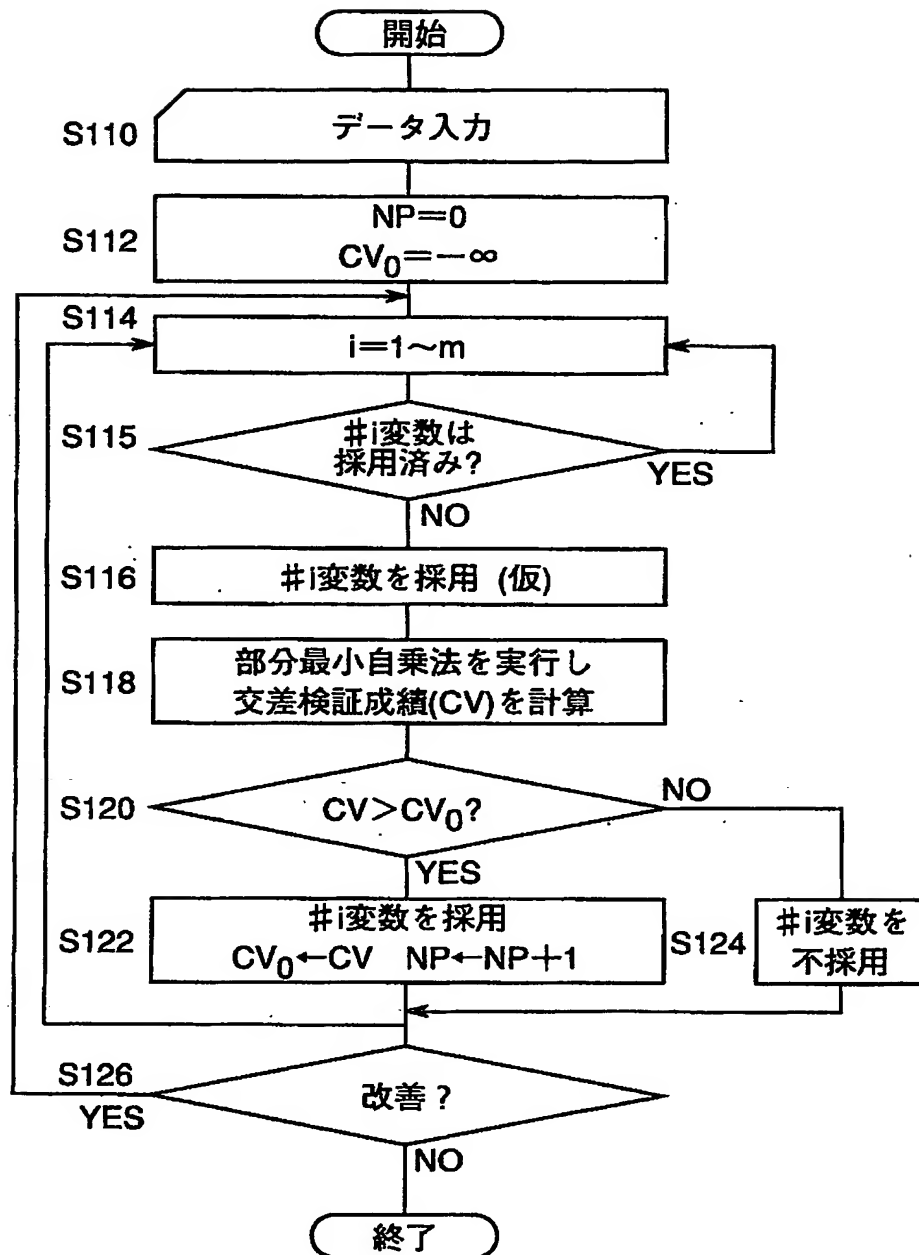
【図 2】



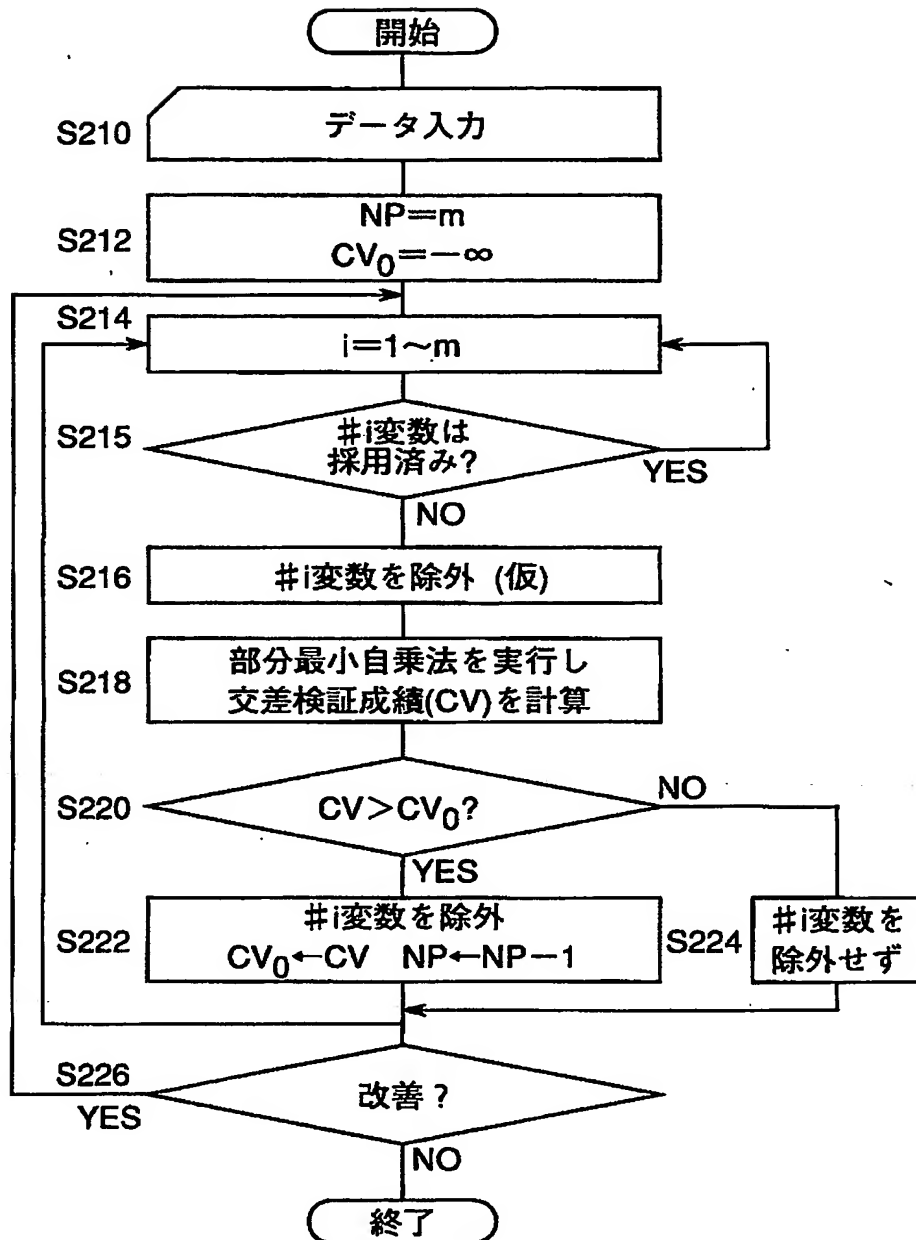
【図 3】



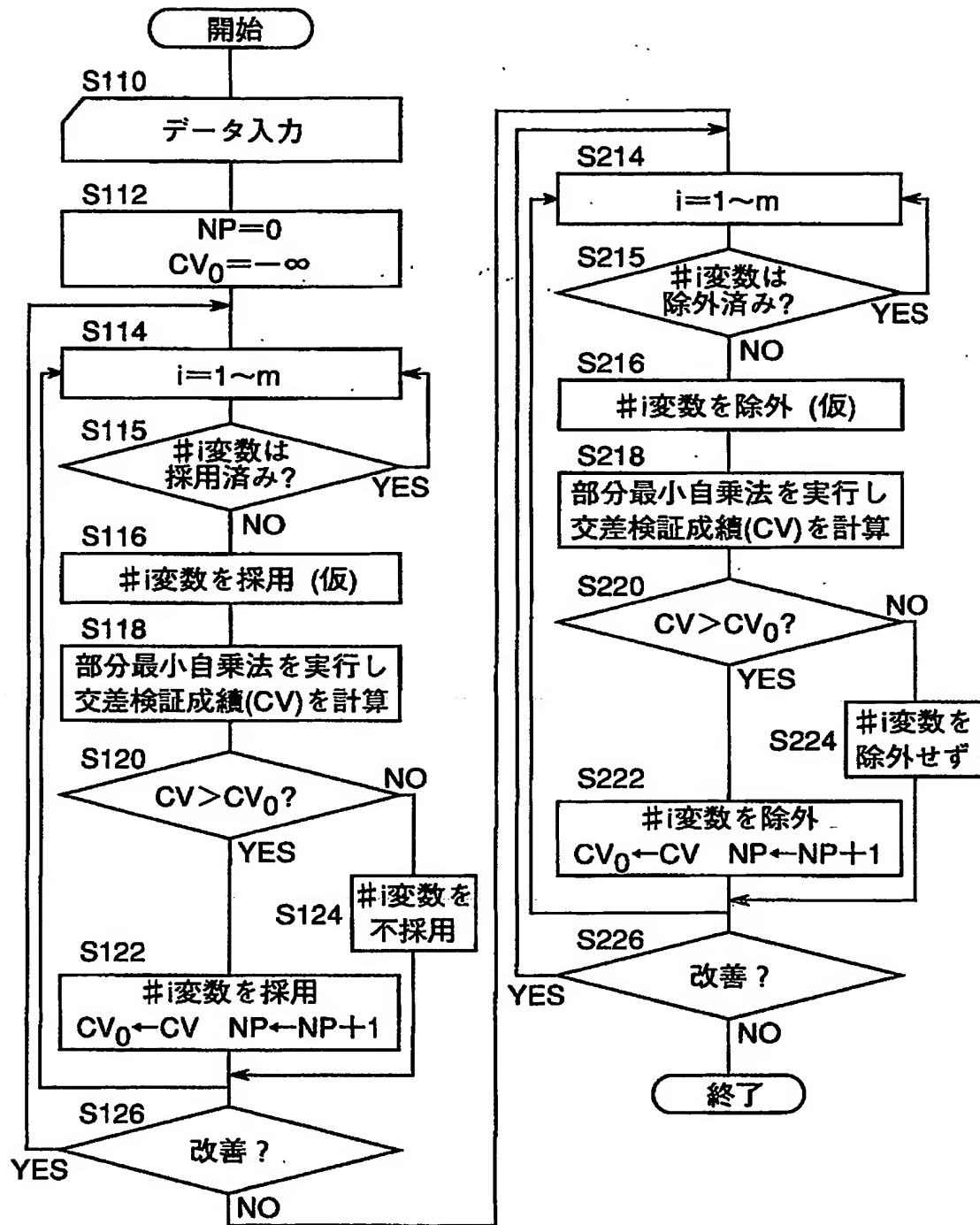
【図 4】



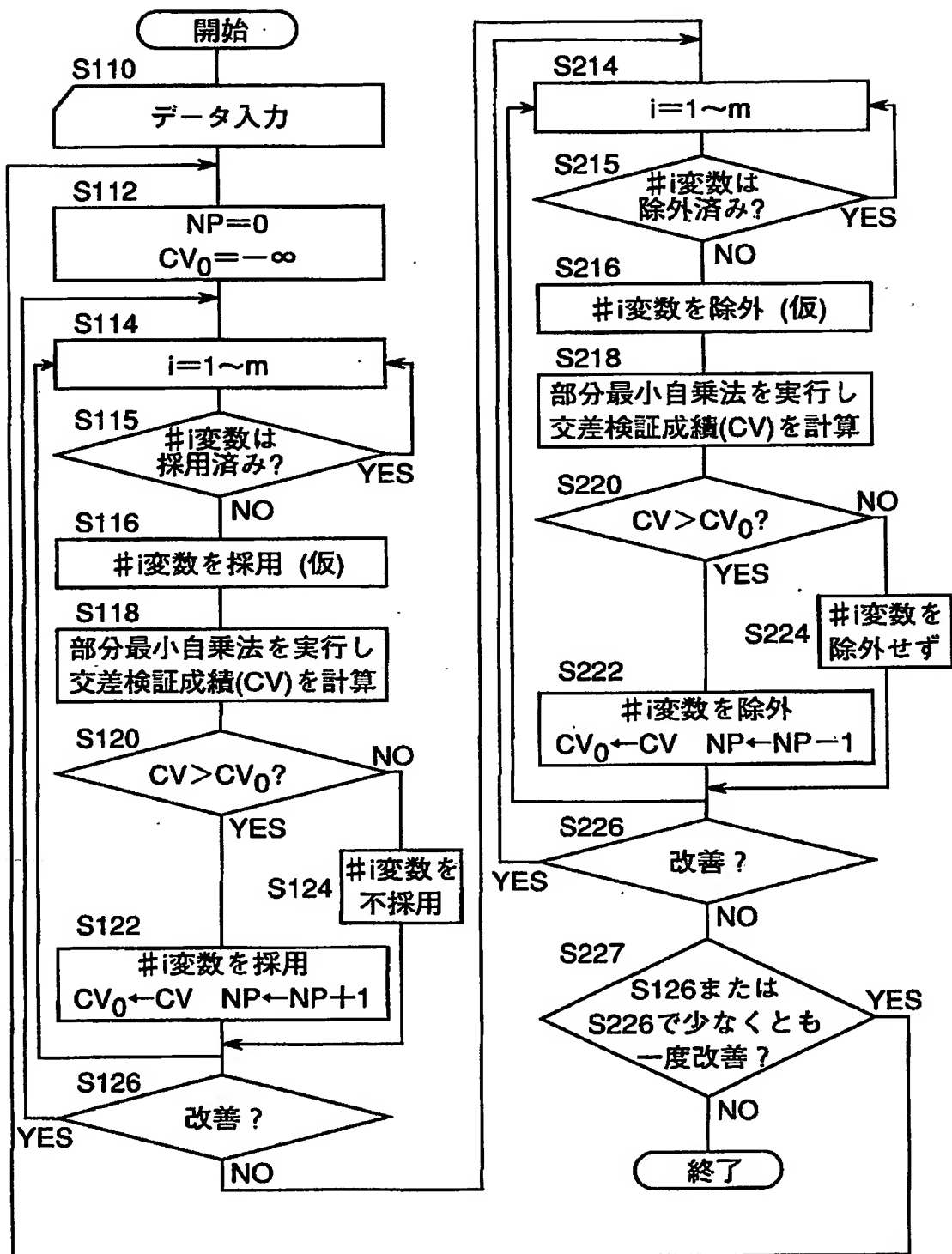
【図 5】



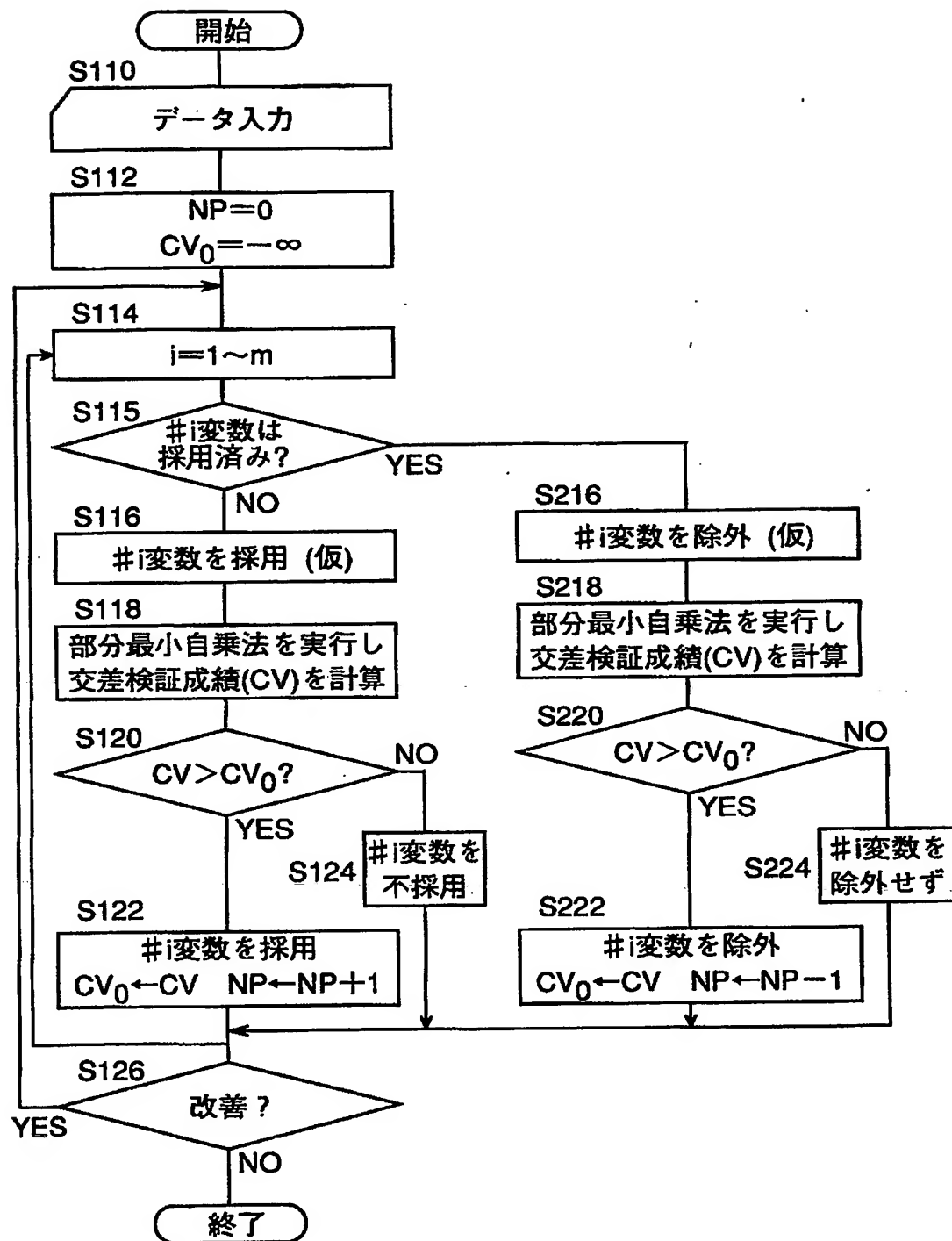
【図 6】



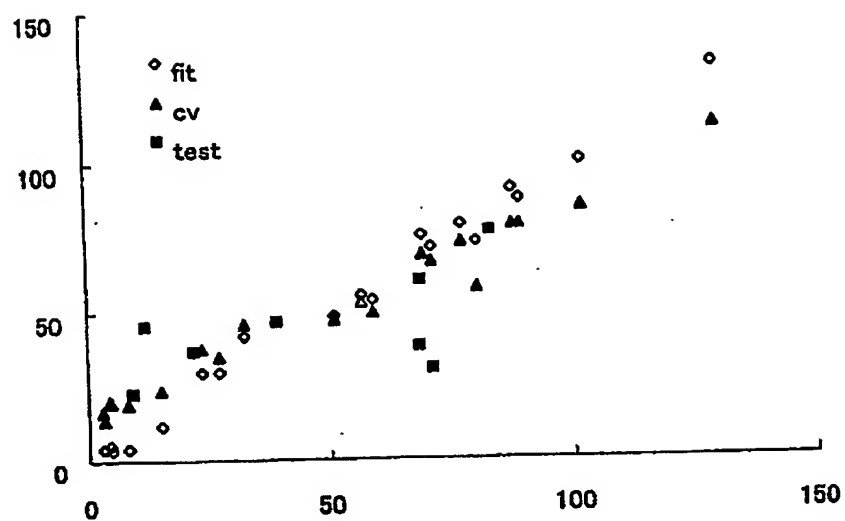
【図 7】



【図 8】



【図9】



【書類名】 要約書

【要約】

【課題】 多変量の遺伝子発現情報の効果的な情報処理を提供する。

【解決手段】 生体の状態と複数の遺伝子発現の量および／または細胞内物質の量との相関モデルを決定するデータ解析において、

生体の状態を目的変数とし、複数(m 個)の遺伝子発現の量および／または細胞内物質の量を説明変数とするデータの集合において、データに含まれる説明変数を選択し、選択された説明変数と目的変数とを含む相関モデルについて交差検証成績を計算し、その結果を評価判定する。ここで、交差検証成績が改善しなくなるまで、説明変数の選択、交差検証成績の計算、その結果の評価判定を行い、部分最小自乗法モデルを決定する。

【選択図】 図2

出 願 人 履 歴 情 報

識別番号 [000000354]

1. 変更年月日 1993年 6月21日

[変更理由] 住所変更

住 所 大阪府大阪市西区江戸堀一丁目3番15号

氏 名 石原産業株式会社